

New Analysis Model Resource Implications

While the new Analysis Model does not explicitly specify resource usage, its requirements clearly have some implications for resource allocation. The largest single resource implication comes through the requirement for access to the Mini. That will require providing adequate disk space for staging a significant fraction of the Mini, which will have some impact on the allocation of disk space for other purposes. Until Mini use patterns are known, it is difficult to estimate how large a fraction of our disk space needs to be allocated for Mini, though the experience of IN2P3 indicates a reasonable value may be 1/2 the size of the Mini. The Mini currently occupies ~10 Tbytes for the full data sample, implying that ~5Tbytes of staging space might be needed.

Since the New Micro format prototype implemented by Eric Charles is actually ~20% *smaller* than the existing Kanga Micro, we have every reason to believe that changing to the New Micro will have at worst a negligible impact on resources. However, during the transition period in which both formats will need to be available, we will clearly need additional resources at the Tier-A sites. We do not however necessarily need to fully replicate both formats at both sites: we could for instance dedicate some Tier-A sites to new format and some to old.

We expect the updated Tag to be significantly smaller than the existing Tag. However, since the tag is already very small, the resource implications for this aspect of the new Analysis Model are negligible.

There are serious resource implications from the Skim requirements of the new Analysis Model. Since skims will be run more frequently, and they will probably be run off the Mini instead of the New Micro, the skim jobs will require additional CPU resources compared to today. Allowing deep copy output will also significantly increase the amount of disk space needed to store the results. Using the 10-series skims as a guide, we estimate that a full set of skims could require between 2.5 and 4 times as much space as a single AllEvents copy for real data, and between 2 and 7 times for generic MC data. The range of values comes from considering different thresholds at which skims would use pointer copies instead of deep copies, which is a policy knob that can be tuned as necessary to meet resource constraints.

The new Analysis Model does have features which somewhat offset the additional resource usage of the new model of skims. For instance, deep copy skims will be more efficient to read than the current pointer skims, which will reduce the CPU and disk IO

used in running analysis jobs. The ability to customize the skim output to include composites should also significantly reduce the CPU used in running analysis jobs. The ability to trim unused content from deep copied events could reduce the per-event size of skimmed data. The fact that skims will be produced more frequently means they will better reflect current needs, which will hopefully encourage more people to run their analysis (more efficiently) on skims instead of on the AllEvents sample, again reducing analysis job CPU usage. The fact that skims in New Micro format could be read either in Framework jobs or directly at the Root prompt implies that some of the space being used for AWG ntuple copies of the micro could be recycled. As these issues all involve significant changes in the use patterns of the skims, it is difficult to estimate their effect quantitatively.